# Exhibit 3

# Smart Reply PRD

go/dynamite-smartreply

**Authors:**
- wenjiazhu@

**Related docs:**
- Eng Design

**Status:** Approved for V1
**Last updated:** August 10, 2018

**Comment [1]:** I have a high level though about PRD for ML features. I found it historically useful to decouple the modeling questions what training architecture, traning data from product definitions. There are two reasons:
 - to make sure product behavior really focuses on what is valuable for users and all the details how this is best surfaced to them and how they interact with the feature
 - based on the above it is easier and potentially quite independent question how to achieve the prediction/ml behavior that fulfills the above requirements. To some extent is best to leave it free of specification. Any model/technilogy that can suppoer the use case is good.

In fact in smart reply gmail we did completely swapped the modeling and training infra over time, actually many times during early development, and even after our original launch, but product behavior stayed stable, up to some small but impactful ui finetunes.

**Comment [2]:** Good points. Tom brought up similar ones as well.  The intention was to have a high-level overview here to make sure that everyone on my team has some basic background. I'll  probably have it link out to a different doc or even remove the section entirely in the future.

# Overview

## Problem

Chat is meant to be a quick and efficient way of communicating with your co-workers. However, users often find themselves struggling to type messages fast enough, or think of the right words to say. Google has built a number of intelligence ML models to help its products be smarter and faster. In particular, various Smart Reply models have demonstrated successful user engagement in Allo and Gmail. The challenge now is to evaluate if one of these ML models can work successfully in Dynamite as well, and further differentiate Dynamite competitively against competitors in the market.

## Goals & Non-Goals

### Goals
*Reinforce Dynamite as the hub for lightweight, frictionless real-time collaboration*

Demonstrate the power of Google intelligence by offering users relevant and interesting Smart Reply suggestions in Dynamite so that they can quickly respond to messages
- Build, evaluate and launch  V1 of Smart Replies in Dynamite with Allo's Smart Suggest API and infrastructure
- Demo at NEXT on July 28
- Achieve ~10% CTR on Smart Reply suggestions to be on par with Gmail & Allo while maintaining positive user satisfaction

### Non-Goals
In the V1, we will avoid touching the current UI/UX of Dynamite and the "Compose" component as much as possible. These other areas are also out of scope:
- Compose box UI/UX overhaul or refresh
- Smart compose - use text content in the Compose box (draft mode) to generate smart replies or responses
- Emoji suggestions or Reaction emoji suggestions, even though emojis may be given as recommendations by the ML model

## Impact

Smart Replies can boost Dynamite's profile as an intelligent chat product and showcase the benefits of adopting products within the Google ecosystem. This feature can also improve user engagement and the speed at which users are able to work productively with their teams.

**Users**

This feature primarily targets <u>Chat end-users</u> who are using Chat as a real-time communication and collaboration product. This feature does not directly target <u>IT Admins</u>, but should meet any criteria they have on accessibility, HR policy, security, etc.

## ML Model Background & Recommendation

For V1, we will use Allo Smart Suggest API to set-up Smart Replies in Dynamite. Based on user feedback and metrics in Dogfood and production, we can then decide if it makes sense to work directly with Kona to build a custom model for Dynamite in a future iteration (V2)

**Allo Smart Suggest**

<u>Allo Smart Suggest</u> is an API layer interface and serving infrastructure powered by ML models owned by 3 different research teams at Google: Hobbes, Kona and Expander. The Smart Suggest team builds a common infrastructure, and incorporates additional business logic on top of the ML models. Smart Suggest also picks between these 3 different ML models based on which one works best for the message. Currently, Smart Suggest's total data training set includes Hangouts Classic, Allo and Google Voice data. They plan to start training on Android SMS data in the near future as well.

| Research Team | ML Model Description |
|---|---|
| <u>Kona</u> | Focused <u>on-server</u> model<br>Trained on Hangouts Classic, Allo and Google Voice data<br>Gmail Smart Reply was built with the Kona model using Gmail data<br><br>**Recommended option for Dynamite** <u>(Details)</u> |
| Hobbes | Focused on <u>on-device</u> model, server-side is not as good<br>Trained on Hangouts Classic and Allo data |
| Expander | <u>On-device</u> model<br>Help facilitate ML teams to get ML on devices<br>*Least relevant option for Dynamite* |

<u>Note:</u> The current SmartReply bot was built using the Smart Suggest API.

**On-Server vs On-Device**

| Type | Pro | Con |
|---|---|---|

| On-Server<br><br>**Recommended**<br><br><u>(Details)</u> | More data, more memory, and thus can provide better recommendations | Privacy concerns (if any)<br><br>User may not want to send their messages to Google. Should NOT be a problem for Dynamite as it is already an on-cloud Google service |
|---|---|---|
| On-Device | Addresses privacy concerns<br><br>Can strip down on-server model from 300MB to 3MB to run on the phone | Smaller dataset, memory restrictions<br><br>E.g. 10,000 English words for Smart Reply on GBoard versus entire English vocabulary |

**Additional Business Logic**

Finally, Allo Smart Suggest provides additional business logic on top of the ML models:

- <u>Language detection:</u> Instead of using the interface language, language is detected based on the last 10 messages sent.
- <u>Diversification</u>: Heuristics to provide a diverse set of reply suggestions instead of showing similar ones.
- <u>Trigger</u>. Use contextual factors to decide whether or not to trigger smart replies, e.g. there is a sensitivity classifier to block smart replies when discussing sensitive topics like tragedy or illness.

# Use Cases & CUJs

## Base Principles

Dynamite's mission is to enable happy productive teams for work. Thus Smart Replies need to respect that mission both in terms of user interaction and content.

- <u>Don't be distracting:</u> Smart replies should not be distracting or disruptive to a user's workflow, either as a message composer or a message recipient.
- <u>Use business appropriate language:</u> All Smart Reply suggestions should be safe and appropriate for work. This means no street slang, curse words or intimate language. There should also be proper capitalization and punctuation.
- <u>User perception of control is important:</u> While we are introducing machine intelligence, users should feel that they have complete control of what they write and share. In other words, "the machine is helping me versus taking over"
- <u>Be a writing aid:</u> Help users feel that they are able to respond *better* and *faster* to messages. We are not necessarily helping users write *more messages to more threads.*
- <u>Preserve primary CUJs:</u> Composing, sending and receiving messages are still the primary workflows. Smart Replies should not compromise the performance or usability of

these key user journeys.

## Use Cases

Today's Model:

- Speed: Think and do less to respond faster
- Agreement: e.g. "Me too" "Can't make it"
- Acknowledgement: e.g. "Got it." "Will do"
- Support:  e.g. "Sounds good." "Let's do it."

Future Applications

- Action suggestions: Integrate with GSuite to recommend docs, meeting times, etc.
- Emoji autosuggest: Suggest emojis for messages or Reactions
- Bot engagement: Learn bot syntax via smart replies for chatting with bots
- Bot discovery: Integrate with Google Assistant to enable new bot discovery
- Customization: Customer-specific language and vocabulary

## CUJs

Below are the primary CUJs for Smart Replies in Rooms, Group DMs and DMs (Video link to animated mocks)

| P | Use Case | Description |
|---|----------|-------------|
| P0 | Displaying smart replies | (1) The user is presented with 3 smart reply options when they click into Reply for any conversation, of which they can select 1.<br>- Rooms: Smart replies are displayed above Compose box only when the user clicks "Reply" to an existing thread or conversation<br>- DMs: Smart replies are *always* displayed above the Compose box<br><br>Allo: Same behavior for DMs<br>Gmail: N/A |
| P0 | Sending smart replies | (2) When the user selects a smart reply option, it is inserted into the Compose box. The user will have an opportunity to either edit the smart reply or add additional content before sending the message.<br>- If we find that most smart replies are sent without any edits, then we can consider sending smart replies automatically when selected<br><br>Allo: Auto sends select reply when tapped<br>Gmail: Same behavior -- allows editing before sending |
| P0 | Ignoring smart replies | (3) The user can choose to ignore any of the 3 smart reply options, and continue to type messages in Compose in the same way they do today<br><br>Allo: Same behavior<br>Gmail: Same behavior |

| P0 | Providing feedback | (4) [Googlers only] The user has the option to provide feedback on any set of Smart Reply options<br>- Call to action for feedback should we prominent<br>- Definitely want to make sure we capture bad suggestions<br>- Mandatory information from user when they submit feedback is the "Bad Suggestion" rating for the quality of the smart replies<br>- Include check-box and disclaimer to allow users to agree to sending data on last 3 messages and the smart replies displayed (auto-checked by default)<br>- Secondary information that would be good to have would be free-text comments, but this step is not required for feedback submission<br>- <u>Feedback flow details</u><br><br><u>Allo:</u> N/A<br><u>Gmail:</u> Same behavior |
|----|----|----|
| P1 | Opting out | (5) The user can opt out of having Smart Replies in the future<br>- Opt out setting is specific to Mobile and Web, eg. user can opt out of SR for Web but still have them on Mobile and vice versa |

# Other Feature Requirements

**Language**

In V1, we will only offer Smart Replies in <u>English</u>. This means that both the interface language *and* the conversation language need to be in English in order for Smart Replies to be displayed. If a conversation language cannot be detected, then we will also not display Smart Replies.

While Allo Smart Reply is currently available in 6 languages (English, French, Spanish, Portuguese, Hindi and Hinglish), curating and vetting a proper whitelist for each will be challenging. Future versions of Smart Reply can either includes these supported languages, or additional languages (e.g. Japanese) via new ML models.

**Whitelist**

The Whitelist is essentially a dictionary of all possible Smart Reply responses. We will start with the current Allo Smart Reply whiteslist, and go through each response to make sure that it's business appropriate:
- Proper capitalization and punctuation
- Remove consumer slang terms, e.g. "That's sick"
- Remove potentially inappropriate phrases such as "I love you" or "honey"

Whitelists are client specific, which means that we can manage and control a specific Whitelist for Smart Replies without affecting other applications.

**Comment [3]:** One consideration to add to this list: even if the whitelist is business friendly how do we assess the risks that the current consumer model might do something not business friendly. E.g. a chat "I heard that we will defrag project X from office Y" suggesting "good news" on this is bad.

Smart reply models on gmail (and I would expect on chat) historically have been rather positively biased. On gmail and I think on allo side as well we had an extra model to detect sensitive topics (tragic events, political topics) and did not trigger on those.

On Gmail side while smart reply is available for work users and luckily the one prominent "non-work friendly" scenario was suggesting "I love you" which we caught in dogfood and later prod feedback didn't raise concerns.

Given that dynamite is solely business focused (to my knowledge), and chat nature is different and models are different we should figure out how much do we care about this "non-business" friendly suggestions and what is worth doing in this direction.

The Dogfood rollout period will be essential to helping us continuously improve the Whitelist via qualitative and anecdotal feedback.

## Other Requirements

Video link to animated mocks

| Req | Description | Timing |
|---|---|---|
| Overflow | We will not allow overflow of smart replies on mobile or web. On web and mobile, if 3 smart replies don't fit on one line, then we'll drop the 3rd one. For mobile specifically, we will allow text wrapping of up to 3 lines within each smart reply pill.<br><br>Since we control our own Whitelist, we can easily assess how often we will be dropping a reply. | Teamfood |
| Existing text | Smart Replies don't get triggered when there is existing text in the Compose box, as in this case the user will already have a reply in mind. This should prevent Smart Replies from triggering when users are actively engaged in a conversation. | Teamfood |
| Overriding smart replies | If a user decides to ignore the reply suggestions and type a response instead, then the Smart Replies will disappear. | Dogfood |
| Last message is from self | Do not trigger Smart Replies if the last message in the DM or Room was sent from yourself<br><br>Note: This case may be covered by the model already | Dogfood |
| Unselected smart replies | Once a Smart Reply suggestion is selected, the other suggestions disappear. | Dogfood |
| Message stream updates - General | [Option 1] Once Smart Replies are displayed for a specific DM or Room Thread (you've already clicked Reply into a thread), no new Smart Replies will appear until you leave the DM for another conversation or become "inactive" on Dynamite and come back.<br><br>[Option 2 - *preferred*] Smart Replies in a DM or Room Thread will update in response to new incoming messages | Teamfood |
| Min. number smart replies | At least 3 smart reply options need to be available, or else none will be displayed | Dogfood |
| History Off | Do not display Smart Replies if History is Off. The assumption is that users often turn History off to discuss sensitive topics. | GA |

**Comment [7]:** Just updated this requirement to having the smart replies disappear once a user starts typing in Compose box.

If we do decide to switch back to having SRs display, then they would still be clickable and be added to your draft in Compose based on your cursor location

**Comment [4]:** +wenjiazhu@google.com, if the user clicks a smart reply at this point, what happens?
_Assigned to Wenjia Zhu_

**Comment [5]:** Just updated this requirement to having the smart replies disappear once a user starts typing in Compose box.

If we do decide to switch back to having SRs display, then they would still be clickable and be added to your draft in Compose based on your cursor location

**Comment [6]:** +wenjiazhu@google.com, if the user clicks a smart reply at this point, what happens?
_Assigned to Wenjia Zhu_

**Comment [8]:** Will start with easiest engineering implementation

| Bots | Bot messages should not be considered as input for generating Smart Reply responses. Therefore, *do not display Smart Replies in DMs with Bots*<br><br>User messages to bots, however, can be considered as input for Smart Replies.<br><br>If the last message in a thread is from a Bot, then do not display Smart Replies (P2) | GA |
|---|---|---|
| Message Edits and Deletion | We should fetch new Smart Replies when messages have been edited or deleted (messages within the last 10) | GA |

## Basic Model Settings

The Allo Smart Suggest team has built a number of "out-of-the-box" model settings that can easily be configured to tune the model without additional engineering effort. For V1, we will keep the default values for each of these settings, and adjust as needed based on the user engagement metrics.

### Data input

The Smart Suggest model will provide recommendations based on messages sent/received on Dynamite. While every message sent on Dynamite will also get sent to Smart Suggest, there is a question on how many messages the model should look at to get smart reply predictions.

| Option | Description |
|---|---|
| Last 10 messages<br><br>**Recommended** | Default option for Allo Smart Suggest<br>ML Model looks at the last 10 messages to come up with smart replies<br><br>*For Rooms, this would need to be limited to a Thread* |
| Last message | Might work better in DMs. E.g. Allo starts providing out-of-context resplies in DMs when looking at last 10 messages. |
| Timestamp | Look at how much time has passed between each message to figure out data input |

> **Comment [9]:** Note: Sergey says that Allo tried models that took time stamp into consideration but didn't have much success

### Number of Smart Replies

You can set how many Smart Replies are presented to users each time (e.g. 2, 3, 4, etc).

For the V1, we will show 3 smart reply options, which aligns with Gmail.

In future iterations, we can experiment with whether showing more smart replies leads to better engagement and quality metrics.

## Content

We can decide which types of Smart Reply messages to prioritize, such as long vs short, emoji versus text.The default version of the model seems to favor short messages, and text over emojis. In the case of emojis, it will still suggest emoji smart replies in specific instances, e.g. you receive an emoji message.

For V1, we will keep the default setting for message length and emoji vs text.

## Trigger rate

The ML model will not always be able to find appropriate Smart Replies -- only a percentage of messages received for each user will come with Smart Reply options. For V1, we will target a ~55% trigger rate to start and adjust as needed.

Note: Allo's trigger rate was about 80%. Android messenger trigger rate is ~50-60%. Gmail smart reply trigger rate is ~30%. There's usually an inverse relationship between Click Through Rate (CTR) and Trigger Rate.

# Work Plan

The first priority is to get a working version of Smart Replies in Dynamite team food or Google Dogfood as soon as possible so that we can evaluate the Smart Suggest model, as well as the user experience of having smart replies business chat.

Gantt Chart

**V1**
- Build V1 of Smart Reply with Smart Suggest API in Dynamite
- Track Metrics in Dogfood
- Tune basic model settings
- Improve the Whitelist
- GA on Web
- English only

**V2**
- Build & GA on mobile
- Revise Whitelist (if needed)
- Work with Kona to customize ML model for Dynamite

- Expand to additional languages

**V3**

- Find a way to train on Dynamite data

# Metrics

We will measure the success of Smart Replies based on metrics from 2 broad categories: Feature engagement and coverage and user satisfaction

- <u>Feature engagement & coverage:</u> Are customers using this feature? How many messages are sent containing a smart reply? How do these metrics compare with Gmail and Allo Smart Reply?
- <u>User satisfaction:</u> Do users like this feature? Do they find Smart Replies distracting? Are Smart Replies appropriate or showing up under the appropriate circumstances? Are we compromising core user flows?

This section primarily covers the former set of quantitative metrics.

## Definitions

Below is how we define our core engagement metrics based on user flow (<u>link)</u>:

> **Comment [10]:** +bal@google.com Hi Balint, I've updated this section with more granular information on how we define the core metrics, and also linked to our doc on events to log (WIP) FYI +tholman@google.com _Assigned to Balint Miklos_
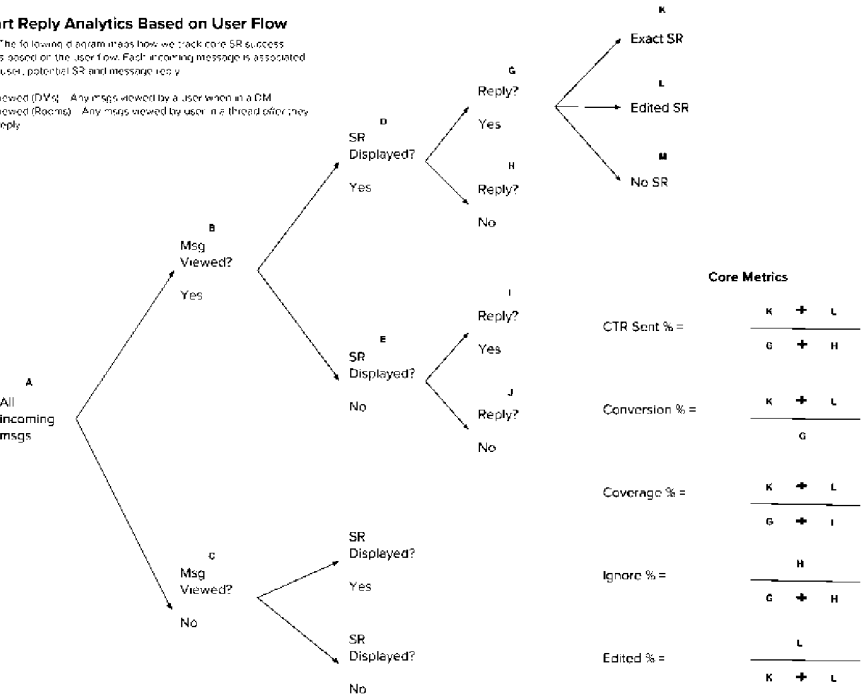
**Smart Reply Analytics Based on User Flow**

Note: The following diagram maps how we track core SR success metrics based on the user flow. Each incoming message is associated with a user, potential SR and message reply.

Msg Viewed (DMs) – Any msgs viewed by a user when in a DM
Msg Viewed (Rooms) – Any msgs viewed by user in a thread after they click Reply

Core Metrics

$$\text{CTR Sent \%} = \frac{K + L}{G + H}$$

$$\text{Conversion \%} = \frac{K + L}{G}$$

$$\text{Coverage \%} = \frac{K + L}{G + I}$$

$$\text{Ignore \%} = \frac{H}{G + H}$$

$$\text{Edited \%} = \frac{L}{K + L}$$

## Full List of Metrics

Thus the following is a comprehensive list of metrics that we would like to track (<u>Metrics Tracker Sheet</u>):

Core metrics

| Type | Metric | Description | Target |
|------|--------|-------------|--------|
| Engagement | CTR (selected)<br><br>- Rooms & DMs | % of Smart Replies displayed (viewed by users) that are selected and inserted into Compose | |
| Engagement | CTR (sent)<br><br>- Rooms & DMs | % of Smart Replies displayed (viewed by users) that are selected *and sent*. E.g. if User A sees Smart Replies 10 times, how often does User A click and send one of the options | TBD |
| Engagement | Conversion rate | Similar to the above, but only looks at cases where the user replied to the message | 10% |

**Comment [11]:** Perhaps a nit, but one thing this does not account for is that users today don't reply to 100% of chat messages that they receive. As such, presumably for a significant % of messages where we show smart reply the user will have no intent to reply (whether manually or via smart reply). Perhaps that's okay, but you potentially want to keep an eye on that and keep in mind baseline of reply frequency in absence of smart reply?

**Comment [12]:** I agree with you, and Balint brought up a similar point in the past as well, which is why we also have the conversion metric, which will only look at cases where the SR was displayed and there was a reply

**Comment [13]:** By a user sees Smart Replies 10 times then send one of the options, does it means when user hover their mouse over an option again and again, each time will be tracked separately? if so we should consider also use Boolean CTR in addition to CTR.

**Comment [14]:** what do you mean by "hover over mouse again and again?" +melissalj@google.com for how we define "SR is seen"

A user clicks into a DM, sees a SR, and doesn't do anything, and then comes back in 5 min, and uses one fo the SRs to respond (nothing has changed) -- I believe we track these separately?

**Comment [15]:** By "hover over mouse again and again?" i mean exactly the same scenario as you mentioned.

| Engagement | Smart replies sent without edits | % of Smart Replies sent without any edits | |
|---|---|---|---|
| Engagement | Trigger rate -<br>-   Rooms & DMs | How often are smart replies displayed when users click into Reply | |
| Engagement | Ignore ratio<br>-   Rooms & DMs | How often smart replies are displayed (viewed by users) but no messages are sent | |
| Reach | Coverage ratio<br>-   Rooms & DMs | What % of all messages sent contain a Smart Reply? | 7% |
| Reach | User penetration | % of unique DAU or MAU sending smart replies | TBD |
| Satisfaction (survey) | Distraction | % of users are "Distracted" or "Very Distracted" by Smart Replies (based on a 5 pt scale) | 25% |
| Satisfaction (survey) | Usefulness | % of users who find Smart Replies "Useful" or "Very Useful" (based on a 5 pt scale) | 75% |
| Satisfaction (survey) | Satisfaction | % of users who find Smart Replies "Satisfying" or "Very Satisfying" | 75% |
| Satisfaction | Opt out rate | Opt out rate for this feature (if available) should be under 10% (TBD, based on Allo criteria for Android messages) | |
| Quality | Active user growth | DAU, WAU and MAU metrics should not drop and growth should be the same as before | |
| Quality | Number of messages sent<br>-   Rooms & DMs | # of messages sent (total and per user) should not drop and stay the same | |
| Performance | Trigger rate - all messages | % of *all incoming messages* received that trigger a smart reply. Allo trigger rate is 80%, Gmail was 30%, Android Messenger is ~60% | 55% |

## Logging

New logging events for measuring Smart Reply metrics are described in this doc.

## All related documents:

- Metrics definitions - user flow
- Metrics tracker
- Metrics event logging

## Risk Mitigation

Via our integration with the Allo Smart Reply Service, we have several risk mitigation measures to quickly prevent Smart Replies from displaying in sensitive situations, or to modify/remove responses from the Whitelist.

Proactively, we have the following methods below. We will be using the Dogfood period to inform us on how we should configure these controls.

- **Sensitivity classification - model**: Allo has an additional sensitivity classification model that helps identify situations where smart replies should *not* be triggered, e.g. messages discussing a tragic event. The sensitivity classifiers are client specific, so they can be customized for Dynamite over time.
- **Sensitivity classification - manual (Trigger Blacklist)**: For messages missed by the sensitivity ML model above, we can also manually specify phrases or words for which smart replies should not be displayed based on regex match, e.g. Don't show smart replies when the word "restructure" or "defrag" is used
- **Blacklist:** Essentially the same as removing a response from the Whitelist, but faster implementation, e.g. Never display "I love you" as a smart reply option.
- **History Off:** Pending confirmation from Dogfood feedback, we are currently planning on suppressing Smart Replies whenever History is Off, as there's a higher likelihood that sensitive information is being discussed in these cases

Post GA, we can also use Allo's Smart Reply Service to do the following in a matter of minutes: (1) Prevent a specific type of response from showing up as a suggestion (Blacklist), (2) Block suggestions when the message content contains specific phrases (Trigger Blacklist)

## Privacy

# Redacted - Privilege

## A11Y & i18N

**Accessibility**

Will need to make sure that users are able to see and/or hear smart reply options, and be able to select a smart reply versus typing in Compose. This could be challenging if the smart reply suggestion is an emoji.

**Comment [16]:** We should see what allo/gmail do for smart replies. SR lets sighted users reply quickly, but in the time non-sighted users listen to the suggestions, it might be faster to just type themselves?

**Comment [17]:** +wenjiazhu@google.com to check with Gmail team

### Internationalization
V1 of Smart Replies will be offered in English only. We will evaluate expansion to other languages in future versions, which may require new Whiteslists and models.

## Feedback & Questions

**Questions**
- What happens when Allo and Hangouts Classic eventually go away in the future and the model becomes driven mostly by Google Voice and Android SMS?
  - Wouldn't worry about Hangout Classic data getting stale as english doesn't really change
  - Would probably need custom model that focuses more on old Hangouts Classic data

# Redacted - Privilege

# Redacted - Privilege

- Determine the difference between different whitelists.

**Squad Feedack**
- How Smart Replies are triggered
  - (Jacob) On mobile, the location of the smart replies on top of the compose box takes away too much real estate -- maybe we can hide the keyboard for DMs and Group DMs
  - (Jacob) Concern that if you start with automatic sending, that people won't go back to it even if you allow editing
  - (Laura) Lots of concern in Gmail around Smart Replies
  - (Senthil) Add latency of ~10ms
  - (Chuan) Two options:
    - Driven by clients. Client makes a server-call given the applicable event - - this could cause a delay -- at least the delay would be human noticeable
    - Pre-cache all smart replies every time a thread or conversation is updated
  - (Eugene + Rahat) Can we always show smart replies if user is @ mentioned or active in it, but for other messages, they get triggered with Compose
  - (Jacob) If I type reply, I already know what I'm going to say.
  - (Senthil) Smart replies are annoying when I know the person very well
- Edit or send automatic
  - (Eugene) prefer edit, because it's consistent with other Compose behavior (e.g. file, photo, etc

- ○ (Senthil)
- (Blockers for Chuan)
    - ○ Is any **human observable delay** acceptable when you click on Reply?
    - ○ Is it important to show different smart replies for each individual in a group thread?
    - ○ Smart replies should not be shown if there is already text in the Compose box